```python
1    from datetime import datetime, timedelta
2    from dateutil.parser import parse
3    import random
4    import pandas as pd
5    import numpy as np
6
7    import sklearn
8    from sklearn.metrics import silhouette_score
9    from sklearn.preprocessing import MinMaxScaler
10
11   import matplotlib.pyplot as plt
12   from matplotlib import cm
13
14   from pyclustering.cluster.kmeans import kmeans
15   from pyclustering.utils.metric import type_metric, distance_metric
16   from pyclustering.utils.metric \
17       import euclidean_distance, manhattan_distance, chebyshev_distance, \
         minkowski_distance
18   from pyclustering.cluster.center_initializer import kmeans_plusplus_initializer
19
20
21   class PycClusteringPeople:
22       df_corona = None  # initial loaded data
23
24       added_column_list = []  # Columns added by calculation within class
25       target_col_name_list = []  # list for target column names
26
27       sse_list = []  # list for displaying SSE of each cluster
28       sil_score_list = []  # list for Silhouette score of clustering result
29       centroids_coord_list = []  # list for storing coordinates of centroids
30
31       # member variable for custom distance function
32       weight_list = []  # weight values list
33
34       cluster_model_dic = {}  # dictionary for model storing
35       scaling_model_dic = {}  # dictionary for model storing
36
37       def __init__(self, file_path, target_col_name_list, base_date):
38           self.df_corona = pd.read_csv(file_path)
39           self.df_corona["Severity"] = self.compute_severity(base_date, self.df_corona)
40           self.added_column_list.append("Severity")
41
42           self.target_col_name_list = target_col_name_list
43
44       def compute_severity(self, base_date, data):
45           """
46           method to preprocess the data for distance function
47           :param base_date: datetime, Base date for calculating elapsed time
48           :return: None
49           """
50           col_num = len(data)  # the number of rows from loaded data
51           severity_list = []  # list for storing severity result
52
53           for i in range(col_num):
54               # selecting specific column to compute 'severity'
55               incur_date_col = data['Incurred Date']
56               status = data['Covid Status']
57               severity = 0  # default is healthy. 0.
58
59               if status[i] == 'Contacted':  # contacted person?
60                   # formula for contacted person:
61                   #   x = 1 - ((today's date) - (infected date)) * 0.05)
62                   elapsed_days = (base_date - parse(incur_date_col[i]).date()).days
63                   severity = (1 - (elapsed_days * 0.05)) * 0.5
64
65               elif status[i] == 'Confirmed':  # confirmed person?
66                   # formula for confirmed person:
```

```python
67              #     x = (1 - ((today's date) - (infected date)) * 0.05)) / 2
68              elapsed_days = (base_date - parse(incur_date_col[i]).date()).days
69              severity = 1 - (elapsed_days * 0.05)
70
71          # add the value to the list
72          # and rounding to solve floating-point problems
73          severity_list.append(round(severity, 4))
74
75      return severity_list
76
77  def display_clustering_result(self,
78                                num_cluster,
79                                cluster_idx_list,
80                                cluster_predicted_list):
81      """
82       function to display clustering result on console as tabular type
83       :param num_cluster: int, the number of cluster
84       :param cluster_idx_list: list, cluster index list, ie. [2, 3, 4, 5, 6]
85       :param cluster_predicted_list: list, result of clustering
86       :return: None
87       """
88
89      severity_list = self.df_corona["Severity"].values.tolist()
90      age_list = self.df_corona["Age"].values.tolist()
91
92      if type(cluster_predicted_list) != list:
93          cluster_predicted_list = cluster_predicted_list.tolist()
94      people_num_of_a_cluster_list = []
95      avg_age_of_a_cluster_list = []
96      avg_severity_of_a_cluster_list = []
97
98      print(f"Number of Clusters: {len(cluster_idx_list)}")
99
100     for cluster_idx in cluster_idx_list:  # 1 cluster
101         num_people = cluster_predicted_list.count(cluster_idx)
102         id_target_data_tuple_list = []
103         target_severity_list = []
104         target_age_list = []
105
106         for person_idx in range(len(cluster_predicted_list)):
107             if cluster_idx == cluster_predicted_list[person_idx]:
108                 target_severity_list.append(severity_list[person_idx])
109                 target_age_list.append(age_list[person_idx])
110                 id_target_data_tuple_list.append((
111                     person_idx + 1,  # [0] of tuple is id
112                     age_list[person_idx],  # [1] of tuple is age
113                     round(severity_list[person_idx], 2)))  # [2] of tuple is
                      severity
114
115         people_num_of_a_cluster_list.append(num_people)
116
117         print(f"\tCluster {cluster_idx}:")
118         print(f"\t\tNumber of People: {num_people}")
119         # print(f"\t\t\t{'ID':<4}{'Age':<4}{'Severity Value'}")
120         # for person_in_cluster in id_target_data_tuple_list:
121         #     print(f"\t\t\t{person_in_cluster[0]:<4}"
122         #           f"{person_in_cluster[1]:<4}"
123         #           f"{person_in_cluster[2]}")
124         print(f"\t\tMinimum of Age values: {min(target_age_list)}")
125         print(f"\t\tMaximum of Age values: {max(target_age_list)}")
126         print(f"\t\tAverage of Age values: "
127               f"{round(sum(target_age_list) / len(id_target_data_tuple_list), 2)}")
128         print(f"\t\tMinimum of Severity values: {min(target_severity_list)}")
129         print(f"\t\tMaximum of Severity values: {max(target_severity_list)}")
130         print(f"\t\tAverage of Severity values: "
131               f"{round(sum(target_severity_list) / len(id_target_data_tuple_list),
                  2)}")
```

```python
132                 print(f"\t\tThe Coordinates of Centroid:")
133                 coords = self.centroids_coord_list[num_cluster - 2][cluster_idx]
134                 print(f"\t\t\tX1 (Severity): {round(coords[0], 2)}")
135                 print(f"\t\t\tX2 (Age): {round(coords[1], 2)}")
136                 try:
137                     avg_age_of_a_cluster_list.append(
138                         round(sum(target_age_list) / len(id_target_data_tuple_list), 2))
139                 except ZeroDivisionError:
140                     avg_age_of_a_cluster_list.append(0)
141
142                 try:
143                     avg_severity_of_a_cluster_list.append(
144                         round(sum(target_severity_list) / len(id_target_data_tuple_list), 2
                        ))
145                 except ZeroDivisionError:
146                     avg_severity_of_a_cluster_list.append(0)
147
148                 print()  # float 1 line
149             self.display_summary_table(people_num_of_a_cluster_list,
150                                        avg_age_of_a_cluster_list,
151                                        avg_severity_of_a_cluster_list)
152         print()  # float 1 line
153
154     def display_load_data(self):
155         """
156          function to display data
157          :return: None
158          """
159         print(f"Total number of People: {len(self.df_corona)}")
160         print(f"{'ID':<4}"
161               f"{'Age':<4}"
162               f"{'Covid Status':<13}"
163               f"{'Severity':<9}"
164               f"{'Address':<10}")
165         for i in range(len(self.df_corona)):
166             print(f"{self.df_corona['ID'][i]:<4}"
167                   f"{self.df_corona['Age'][i]:<4}"
168                   f"{self.df_corona['Covid Status'][i]:<13}"
169                   f"{round(self.df_corona['Severity'][i], 3):<9}"
170                   f"{self.df_corona['Address'][i].split()[0]:<10}"
171                   )
172         print()  # float 1 line
173         grouped_status = self.df_corona['Severity'].groupby(self.df_corona['Covid
              Status'])
174
175         print(f"Number of healthy people: {grouped_status.count()['Healthy']}")
176         print(f"Number of contacted people: {grouped_status.count()['Contacted']}")
177         print(f"Number of confirmed people: {grouped_status.count()['Confirmed']}")
178
179         print(f"Average Severity of contacted people: "
180               f"{round(grouped_status.mean()['Contacted'], 2)}")
181         print(f"Average Severity of confirmed people: "
182               f"{round(grouped_status.mean()['Confirmed'], 2)}")
183         print()  # float 1 line
184
185     def display_summary_table(self,
186                               people_num_of_a_cluster_list,
187                               avg_age_of_cluster_list,
188                               avg_severity_of_cluster_list):
189         """
190          function to display the data as tabular summary
191          :param people_num_of_a_cluster_list: list, the number of people in a cluster
192          :param avg_age_of_cluster_list: list,
193          :param avg_severity_of_cluster_list:
194          :return:
195          """
196         len_id = 17
```

```
197            len_p_num = 11
198            len_age = 13
199            len_sev = 15
200            len_sum = len_id + len_p_num + len_age + len_sev
201
202            # top row
203            print(f"\t{'-' * (len_sum + 11)}")
204            print(f"\t{'Cluster ID':>{len_id}} "
205                  f"| {'# of People':>{len_p_num}} "
206                  f"| {'Avg. of Ages':>{len_age}} "
207                  f"| {'Avg. of Severity':>{len_sev}} ")
208
209            # contents of table
210            cluster_id = 0
211            for people_num, avg_age, avg_sev in zip(people_num_of_a_cluster_list,
212                                                    avg_age_of_cluster_list,
213                                                    avg_severity_of_cluster_list):
214                print(f"\t{cluster_id:>{len_id}} "
215                      f"| {people_num:>{len_p_num}} "
216                      f"| {avg_age:>{len_age}} "
217                      f"|{avg_sev:>{len_sev}}")
218                cluster_id += 1
219
220            print(f"\t{'-' * (len_id + 1)}"
221                  f"|{'-' * (len_p_num + 2)}"
222                  f"|{'-' * (len_age + 2)}"
223                  f"|{'-' * (len_sev + 2)}-")
224
225            # bottom row
226            print(f"\t{'Total':^{len_id}} | {sum(people_num_of_a_cluster_list):>{len_p_num
                  }} |")
227            print(f"\t{'SSE':^{len_id}} | {round(self.sse_list[len(
                  people_num_of_a_cluster_list) - 2], 2):>{len_p_num}} |")
228            print(f"\t{'Silhouette Score':>{len_id}} "
229                  f"| {round(self.sil_score_list[len(people_num_of_a_cluster_list) - 2], 2
                  ):>{len_p_num}} |")
230            print(f"\t{'-' * (len_sum + 11)}")
231
232     def draw_graph(self):
233         """
234          method to draw clustering result
235          :return: None
236         """
237         pass
238
239     def draw_silhouette(self):
240         """
241          method to draw graph using silhouette scores
242          :return: None
243         """
244         pass
245
246     def draw_elbow_method(self, sse_list):
247         """
248          method to draw elbow graph using SSE(Sum of Squares Error)
249          :param sse_list: list of SSE
250          :return: None
251         """
252         plt.plot(range(2, 10), sse_list, marker='o')
253         plt.xlabel("The Number of Cluster")
254         plt.ylabel("SSE")
255         plt.show()
256
257     def describe_id(self, new_point, cluster_num, cluster_model):
258         # TODO: Display the information of reason that why the data in in the cluster.
259         #   id is cluster ID. then, this function should explain about the closest
             cluster.
```

```python
260            #    - 경계선인가? 그렇다면 퍼센티지로 나타낼 수 있는가?
261            #    - plotting
262            # Done
263            #    - 가장 가까운 클러스터의 거리와 두번째로 가까운 클러스터의거리
264            #    - 어떤 원소들이 여기에 속하는가?
265            centroid_coords_list = self.centroids_coord_list[cluster_num-2]
266
267        idx_distance_tuple_list = []
268        for idx in range(len(centroid_coords_list)):
269            idx_distance_tuple_list.append(
270                (idx, self.weighted_euclidean_distance(new_point, centroid_coords_list[
271                    idx])))
272
273        print("\tDistance List (Top 3 nearest Clusters)")
274        sorted_list = sorted(idx_distance_tuple_list, key=lambda x: x[1])
275        closest_cluster_id = sorted_list[0][0]
276
277        print(f"\t\t{'Cluster ID':>11} |{'Distance':>9}")
278        for idx_distance_tuple in sorted_list:
279            if sorted_list.index(idx_distance_tuple) > 2:
280                break
281            print(f"\t\t{idx_distance_tuple[0]:>11} |{round(idx_distance_tuple[1], 3)
282                :>9.3f}")
283        print()  # float 1 line
284
285        feature_values = self.df_corona[self.target_col_name_list]
286
287        # display part
288        print(f"\tList of data belonging to the cluster {closest_cluster_id}:")
289
290        closest_cluster_elements_list = cluster_model.get_clusters()[
291            closest_cluster_id]
292        for data_id in closest_cluster_elements_list:
293            if closest_cluster_elements_list.index(data_id) % 5 == 0:
294                print(f"\t\t{str(feature_values.iloc[data_id, :].values.tolist())
295                    :<14}", end='')
296            elif closest_cluster_elements_list.index(data_id) % 5 == 4:
297                print(f"{str(feature_values.iloc[data_id, :].values.tolist()):<14}")
298            else:
299                if closest_cluster_elements_list.index(data_id) + 1 == len(
300                    closest_cluster_elements_list):
301                    print(f"{str(feature_values.iloc[data_id, :].values.tolist()):<14}")
302                else:
303                    print(f"{str(feature_values.iloc[data_id, :].values.tolist()):<14}"
304                        , end='')
305        print()  # float 1 line
306
307    def find_cluster(self, new_person_data, model, base_date=datetime.today().date()):
308        # preprocess of input data
309        new_person_data["Severity"] = self.compute_severity(base_date, new_person_data)
310
311        target_data = new_person_data[self.target_col_name_list].__deepcopy__()
312        # scaling some columns
313        scale_col_name = ['Age']
314        for col_name in scale_col_name:
315            target_data = self.scale_column(target_data, col_name,
316                                            using_enrolled_model=True)
317
318        new_point = target_data.loc[:0, tuple(self.target_col_name_list)].values.
319            tolist()
320        # predict result: [cluster_id, cluster_id, ... ,]
321
322        return model.predict(new_point)[0], new_point[0], new_person_data
323
324    def get_cluster_model_dic(self):
325        """
326         function to return cluster_model_dic
```

```python
320          :return: list, cluster_model_dic
321          """
322          return self.cluster_model_dic
323
324      def get_scaling_model_list(self):
325          """
326          function to return scaling_model_dic
327          :return: list, scaling_model_dic
328          """
329          return self.scaling_model_dic
330
331      def initialize_random_centroid(self, num_centroid, is_0_1_normalized=True):
332          """
333          A function that generates a centroid of random coordinates-
334          -as many as the number of clusters received.
335          :param num_centroid: int, the number of centroids
336          :param is_0_1_normalized: boolean, Whether the feature is normalized
337          :return: the random coordinates of centroids list
338          """
339          if is_0_1_normalized:  # are all columns normalized to 0-1?
340              # for i in num_centroid:
341              return [[random.uniform(0, 1), random.uniform(0, 1)] for _ in range(
                     num_centroid)]
342
343          else:  # there is an unnormalized column
344              pass
345
346      def pyc_cluster_kmeans(self,
347                             num_cluster,
348                             weight_list,
349                             distance_function):
350          """
351          function to cluster data
352          :param num_cluster: int, the number of clusters
353          :param weight_list: list, weight list of features
354          :param distance_function: string, the abbreviation of distance function
355          :return: list, clustered result
356          """
357          self.weight_list = weight_list
358
359          # my_distance_function = lambda p1, p2: p1[0] + p2[0] + 2
360          if distance_function == 'eu':
361              # metric = distance_metric(type_metric.EUCLIDEAN)
362              metric = euclidean_distance
363          elif distance_function == 'ma':
364              # metric = distance_metric(type_metric.MANHATTAN)
365              metric = manhattan_distance
366          elif distance_function == 'mi':
367              # metric = distance_metric(type_metric.MINKOWSKI)
368              metric = minkowski_distance
369          elif distance_function == 'c_eu':
370              metric = distance_metric(type_metric.USER_DEFINED, func=self.
                     weighted_euclidean_distance)
371
372          target_data = self.df_corona.loc[:, self.target_col_name_list]  # To select
                 required data
373
374          scale_col_name = ['Age']
375          for col_name in scale_col_name:
376              target_data = self.scale_column(target_data, col_name)
377          # if "Age" in self.target_col_name_list:  # scaling only "Age" column
378          #     age = target_data.loc[:, "Age"].values.reshape(-1, 1)
379          #     scaler = preprocessing.MinMaxScaler()
380          #     # data = scaler.fit_transform(data)   # To scale data from 0 to 1
381          #     target_data.loc[:, "Age"] = scaler.fit_transform(age)
382
383          # set the number of data and centroids
```

```python
384             self.data_cal_count = num_cluster * len(target_data)
385             self.num_of_data = num_cluster * len(target_data)
386             self.cent_cal_count = 1
387             self.num_of_cent = 1
388
389             # initializing centroids
390             # initial_centers = kmeans_plusplus_initializer(target_data,
                num_cluster).initialize()
391             # initial_centers = self.initialize_random_centroid(num_cluster)
392             initial_centers = [[i * 0.1, i * 0.1] for i in range(num_cluster)]
393
394             kmeans_instance = kmeans(target_data, initial_centers, metric=metric)
395             self.cluster_model_dic[num_cluster] = kmeans_instance
396
397             kmeans_instance.process()
398             clustered_list = kmeans_instance.get_clusters()
399             clustered_list = self.pyc_result_to_column(clustered_list, len(target_data))
400
401             # add the column
402             self.df_corona['Cluster ID'] = clustered_list
403
404             # storing the coordinates of centroids
405             self.centroids_coord_list.append(kmeans_instance.get_centers())
406
407             # storing SSE(Sum of Squared Errors)
408             self.sse_list.append(kmeans_instance.get_total_wce())
409
410             # strong Silhouette Score
411             self.sil_score_list.append(silhouette_score(target_data, clustered_list))
412
413             return clustered_list
414
415     def pyc_result_to_column(self, pyc_cluster_result, people_num):
416         """
417          function to change shape of Pyclustering to pandas
418          :param pyc_cluster_result: nd list, result of clustering using Pyclusterin
419          :param people_num: int, the number of people(data)
420          :return: list, re-shaped list
421         """
422         clustered_list = [0 for _ in range(people_num)]
423
424         cluster_id = 0
425         for id_list_of_a_cluster in pyc_cluster_result:
426             for idx in id_list_of_a_cluster:
427                 clustered_list[idx] = cluster_id
428             cluster_id += 1
429
430         return clustered_list
431
432     def plot_data(self, num_cluster, additional_data=None):
433         """
434          To plot result
435          :return:
436         """
437         groups = self.df_corona.groupby("Cluster ID")
438         fig, ax = plt.subplots()
439         for name, group in groups:
440             ax.plot(group.Severity, group.Age, marker='o', linestyle="", label=name)
441         if additional_data is not None:
442             ax.plot(additional_data[0], additional_data[1],
443                     marker='*', linestyle="", label='New Data', markersize=15)
444
445         ax.legend(fontsize=12)
446         plt.title("Result of Clustering (K=" + str(num_cluster) + ", Weight=" + str(
                self.weight_list) + ")")
447         plt.xlabel("Severity")
448         plt.ylabel("Age")
```

```python
449            # plt.show()
450            file_path = "./Cluster_Result_Plotting_pyc/cluster_result_" + str(num_cluster)
                 + "_" + str(self.weight_list) + ".png"
451            if additional_data is not None:
452                file_path = file_path[:-4] + '_new_data_plot_.png'
453            fig.savefig(file_path, dpi=300)
454            plt.close()
455
456        def scale_column(self, data_frame, col_name, feature_range=(0, 1),
457                         model_enroll=True,
458                         using_enrolled_model=False):
459            """
460             function to scale data as 0~1
461             :param data_frame: pandas dataframe, original data.
462             :param col_name: string, column name to scale
463             :return: scaled data frame
464            """
465            origin_data = data_frame.loc[:, col_name].values.reshape(-1, 1)
466            if using_enrolled_model:
467                model_enroll = False
468                scaler = self.scaling_model_dic[col_name]
469                data_frame.loc[:, col_name] = scaler.transform(origin_data)
470            else:
471                scaler = MinMaxScaler(feature_range=feature_range)
472                data_frame.loc[:, col_name] = scaler.fit_transform(origin_data)
473
474            if model_enroll:  # storing the scaling model
475                self.scaling_model_dic[col_name] = scaler
476
477            return data_frame
478
479        def save_as_csv(self, num_cluster):
480            """
481             function to save data as .csv file
482             :param num_cluster: int, the number of cluster.
483             :return: None
484            """
485            temp_df = self.df_corona.__deepcopy__()
486
487            file_name = f"clustered_corona_data_k={num_cluster}_" \
488                        f"{'Severity_Age'}_{''.join(str(self.weight_list))}.csv"
489            temp_df.to_csv(file_name, encoding='utf-8-sig')
490
491        def weighted_euclidean_distance(self, point1, point2):
492            """
493             custom distance function
494             :param point1: list, list of feature values or coordinates list of centroid
495             :param point2: coordinates list of centroid
496             :return: distance between point1 and point2
497            """
498            distance = 0.0  # distance between point1 and point2
499
500            # when calculating the distance of two coordinates
501            if np.shape(point1) == (2,):
502                # point 1 is data.
503                # point 2 is centroid
504                for weight, p1_coord, p2_coord in zip(self.weight_list, point1, point2):
505                    distance += weight * (p1_coord - p2_coord) ** 2.0
506
507            else:  # when updating Centroid
508                # point 1 and 2 are centroid
509                for prev_cent, curr_cent in zip(point1, point2):
510                    for weight, pc_coord, cc_coord in zip(self.weight_list, prev_cent,
511                        curr_cent):
512                        distance += weight * (pc_coord - cc_coord) ** 2.0
513
514            return distance ** 0.5
```